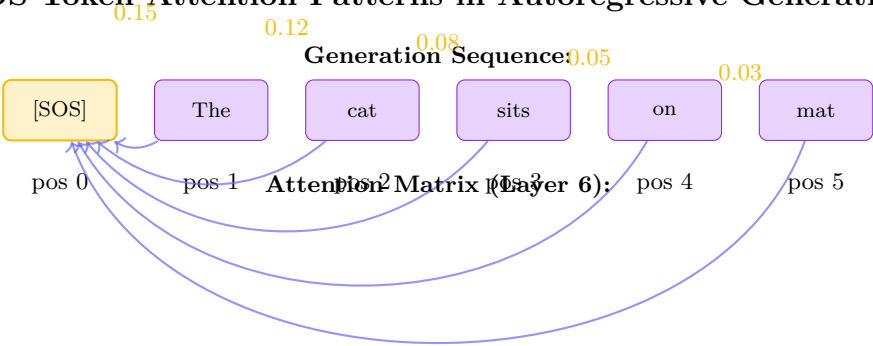


SOS Token Attention Patterns in Autoregressive Generation



Key Insights

[SOS]

Token A

Persistent conditioning signal

Attends to SOS for global context

Bidirectional attention flow

Causal masking prevents future information leakage

Generation Process: Each token uses SOS as initialization context

